

Navigating and Analyzing Large Economic Datasets_{part2}

Xiaoyi Zhao

Virginia Tech Data Science for the Public Good Program

xiaoyiz@vt.edu

Table of Contents

- 1 How to Clean the dataset
- 2 How to organize your code

Checklist for Thoughtful Cleaning

Before You Clean, Ask Yourself...

- What am I trying to measure?
- What assumptions am I making when I clean my data?

What are the things you need to look into when we clean

- Units and scale:
 - For example, is income measured in dollars, thousands, or percentiles?
- Model structure matters:
 - Are you planning to use a linear model? Make sure your variable transformations align with that choice.

Tips

- Name the variables descriptively so someone else (or future you) can follow your work.
 - good example: age_female
 - bad example: age_2
- Keep track of how you transform raw variables (e.g., categorizing income levels or recoding region)

Tips

- Being consistent with file name and variable names
 - For example, for the data acs_2000, set the name acs_2000_i(after clean the allocated), acs_2000_pre(after generate the variables), and acs_2000_final
- Validate at each step
- Separate your data cleaning steps into modular files, and manage the full workflow through a master script—so anyone reviewing your work only needs to run that one file.

Thank You for listening!

Questions and Discussion.